# CBT-BENCH: Evaluating Large Language Models on Assisting Cognitive Behavior Therapy

**Mian Zhang**[*1], **Xianjun Yang**[*2], **Xinlu Zhang**[2], **Travis Labrum**[3],
**Jamie C. Chiu**[4], **Shaun M. Eack**[3], **Fei Fang**[5], **William Yang Wang**[2], **Zhiyu Zoey Chen**[1]

[1]Department of Computer Science, The University of Texas at Dallas,
[2]Department of Computer Science, University of California, Santa Barbara,
[3]School of Social Work, University of Pittsburgh,
[4]Department of Psychology, Princeton University,
[5]School of Computer Science, Carnegie Mellon University
{mian.zhang, zhiyu.chen2}@utdallas.edu, xianjunyang@cs.ucsb.edu

## Abstract

There is a significant gap between patient needs and available mental health support today. In this paper, we aim to thoroughly examine the potential of using Large Language Models (LLMs) to assist professional psychotherapy. To this end, we propose a new benchmark, CBT-BENCH, for the systematic evaluation of cognitive behavioral therapy (CBT) assistance. We include three levels of tasks in CBT-BENCH: **I: Basic CBT knowledge acquisition**, with the task of multiple-choice questions; **II: Cognitive model understanding**, with the tasks of cognitive distortion classification, primary core belief classification, and fine-grained core belief classification; **III: Therapeutic response generation**, with the task of generating responses to patient speech in CBT therapy sessions. These tasks encompass key aspects of CBT that could potentially be enhanced through AI assistance, while also outlining a hierarchy of capability requirements, ranging from basic knowledge recitation to engaging in real therapeutic conversations. We evaluated representative LLMs on our benchmark. Experimental results indicate that while LLMs perform well in reciting CBT knowledge, they fall short in complex real-world scenarios requiring deep analysis of patients' cognitive structures and generating effective responses, suggesting potential future work.[1]

## 1 Introduction

Mental health conditions have reached alarming levels globally, with one in eight people affected, according to the World Health Organization (2023)[2]. There is a severe shortage of mental health professionals. In the U.S., more than 160 million people live in areas with insufficient mental health providers, with rural regions being especially underserved[3]. This critical gap underscores the need for AI-driven tools to support professionals and expand access to care. Existing research has explored mental health condition classifications (Gao et al., 2018; Senn et al., 2022), empathetic conversations (Sharma et al., 2021, 2023a; Adikari et al., 2022), and chatbots designed for simple discourse structures (Hsu et al., 2023). However, work on professional assistance in real therapy settings remain limited. Some studies have addressed specific tasks in cognitive behavioral therapy (CBT) (Beck, 2020), such as cognitive distortion classification (Chen et al., 2023b). However, there are many other critical stages in CBT that could potentially be enhanced through automation.

In this work, we aim to thoroughly investigate the proficiency and potential of LLMs in supporting various facets and stages of professional mental health care. To this end, we propose CBT-BENCH, a systematic benchmark for evaluating CBT efficacy. CBT-BENCH is structured in three levels, from CBT knowledge recitation to the therapeutic responses generation in CBT sessions, providing a hierarchical assessment of CBT capabilities. To ensure the professionalism and high quality of our benchmark, we collaborate with domain experts (clinical psychologists, professors, and social workers) throughout the construction of CBT-BENCH.

In **level I**, we aim to assess **basic CBT knowledge acquisition**. We propose CBT-QA, a new dataset of 220 multiple-choice questions. The QA pairs are collected from CBT exam questions for Master of Social Work (MSW) students as well as compositions from CBT experts, covering a wide range of CBT knowledge, including basic concepts, practical knowledge, and case studies.

In **level II**, we aim to assess **cognitive model**

---

[3]https://bhw.hrsa.gov/data-research/review-health-workforce-research

**understanding**. Modeling how the cognitive components of the patients, such as the beliefs and thoughts, are connected, is at the core of CBT (Beck, 2020; Kuyken et al., 2011). In this work, we propose CBT-CD, a new dataset of cognitive distortion classification with 146 high-quality examples. We also propose two new tasks to assess the understanding of patients' beliefs: primary core belief classification and fine-grained core belief classification. For primary core belief classification, we construct CBT-PC, a new dataset of 184 examples with three primary core belief categories. For fine-grained core belief classification, we propose CBT-FC, a new dataset of 112 examples with nineteen fine-grained core belief categories. Building AI models for such tasks to assist the cognitive modeling process has the great potential to enhance therapists' accuracy and productivity.

In **level III**, we aim to evaluate **therapeutic response generation**, the ultimate CBT efficacy – whether the model can effectively respond to patient speech during CBT sessions like the therapists. Due to the privacy constraints associated with real CBT session data, accessing extensive datasets is challenging. In collaboration with professors specialized in CBT education, we propose CBT-DP, using *Deliberate Practice* as outlined in Boswell and Constantino (2022) for our assessment. This methodology is traditionally employed in assessing CBT proficiency among graduate students. It encompasses a collection of exercises to respond to typical patient speeches across ten key aspects of CBT sessions, categorized into three levels of difficulty, totaling 156 distinct exercises. This approach ensures comprehensive coverage of critical and challenging scenarios likely to occur in real sessions. Generating high-quality responses to patient speeches in CBT-DP can serve as an effective proxy for real CBT session efficacy.

We experiment with six popular LLMs for level I and level II in CBT-BENCH and find that 1) the models of large sizes are better at answering CBT knowledge questions; 2) simply making the models larger could not enhance their understanding ability of the cognitive model; 3) current LLMs struggle with detecting fine-grained cognitive disorders or core beliefs. A deeper analysis of the model performance for each class points out the main reasons for a wrong prediction and directions to improve the corresponding ability of the models. For the level III task, we find that LLMs generally follow a rigid logical reasoning process but lack a crucial

skill in psychotherapy—thinking and guiding from the patient's perspective to respect their autonomy and build rapport. These findings highlight significant limitations in applying LLMs to real-world psychotherapy practice and provide valuable insights for future research and development toward more accessible and efficient mental health care.

## 2 Related Work

Our work is the first to systematically evaluate LLMs' ability to assist professional human therapists in the specialized field of CBT. The works directly related to our research include injecting domain knowledge of mental health into the models (Yang et al., 2024; Kim et al., 2024), cognitive disorder detection (Shreevastava and Foltz, 2021a; Wang et al., 2023; Chen et al., 2023b), negative thoughts recognition and reframing (Sharma et al., 2023b; Maddela et al., 2023; Sharma et al., 2024), and patient simulation (Chen et al., 2023a) or therapist simulation (Liu et al., 2023) in therapeutic conversation. Wang et al. (2024) proved that with clear modeling of the cognitive model of patients, LLMs could act more like real patients. Also, we believe that enhancing the understanding of the cognitive models of patients could be beneficial to interpretable medical decisions (Yang et al., 2023), which is a crucial step towards reliable and safe automated mental healthcare (Ji et al., 2023; Grabb et al., 2024). In the proposed therapeutic response generation task, we leverage the feedback from professional therapists to assess the potential of LLMs responding like a therapist in real therapy sessions, while Louie et al. (2024) explored how to make LLMs roleplay patients with the feedback from professional therapists.

## 3 CBT-BENCH

In this section, we elaborate on how CBT-BENCH is constructed. We discuss **level I, II, and III** tasks in §3.1, §3.2, and §3.3, respectively. Note that apart from collaborating with social work professors and clinical psychologists, who are our co-authors, we recruit CBT experts (clinical psychologists, social workers, etc.) from UpWork[4] for all additional data annotation tasks (See §8 for recruiting details).

| Knowledge Types | Example QA Pairs from CBT-QA | Distributions (%) |
|---|---|---|
| Basic CBT knowledge and concepts | Albert Ellis' Cognitive Model includes which components?<br>**A.** Activating Events – Behaviors – Cognitions, **B.** Antecedents – Beliefs – Consequences,<br>**C:** Activating Events – Beliefs – Consequences, **C.** Antecedents – Behaviors – Consequences | 41.82 |
| Practical CBT knowledge | When helping clients evaluate automatic thoughts,<br>therapists should generally help clients evaluate which aspects of those thoughts?<br>**A.** Accuracy and/or intensity, **B.** Intensity and/or utility,<br>**C.** Accuracy and/or utility | 34.09 |
| Case studies | The client has identified an automatic thought of "My partner is going to break up with me".<br>The therapist asks the client, "If your thought is accurate and your partner does break up with you,<br>what does that mean about you?". The therapist is most likely trying to identify:<br>**A.** The client's intermediate belief, **B.** The client's core belief,<br>**C.** The client's thinking error | 18.18 |
| Others | What are some ways that CBT therapists can engage in therapy from a multicultural perspective?<br>(select all that apply)<br>**A.** Not take clients from a different culture than their own,<br>**B.** Ask clients about the strengths and challenges of their cultural, racial, and ethnic identity during intake,<br>**C.** Being aware of their own cultural values and biases,<br>**D.** Work together with the client to incorporate the client's core values | 5.9 |

Table 1: Knowledge types in CBT-QA, with example QA pairs and distributions in the test set.

## 3.1 Level I: Basic CBT Knowledge Acquisition

To evaluate the basic CBT knowledge acquired by LLMs, we propose a new dataset, CBT-QA, encompassing 220 multiple-choice questions. We first worked with two social work professors and collected 92 multiple-choice questions from the exams they issued for graduate CBT courses. Each question has 2-5 answer choices. Then, we hired four CBT experts to compose new QA pairs, using the QA pairs from the exam source as guidance, and each expert was tasked with composing 50. To ensure high quality, we required the experts to cross-verify the QA pairs and excluded those deemed of low quality. We ended up with 178 high-quality QA pairs, which, when combined with the 92 ones sourced from exams, resulted in our CBT-QA dataset of 270 QA pairs. We randomly selected 50 ones for use such as in-context learning examples. The remaining 220 pairs were designated as the final test set.

To get a better understanding of the fine-grained knowledge types involved in CBT-QA, we employed another two CBT experts to categorize the test set into four categories. Table 1 shows example QA pairs from the four categories and corresponding distributions in the test set. See Appendix B for more examples. To assess human performance, we employed another two CBT experts to solve the test set and ended up with an average accuracy of 90.7% with an agreement rate over 80%.

## 3.2 Level II: Cognitive Model Understanding

In the exam setting of CBT-QA, most questions assess the models' ability to recite knowledge, as

LLMs may have seen this information during pre-training—especially when it comes to basic CBT knowledge and concepts that are prevalent online. However, knowledge recitation alone is insufficient to support the effective use of LLMs in assisting real practice. To address this, we propose the Level II task set: cognitive model understanding. In CBT, therapists develop a cognitive model to represent a patient's unhealthy cognitive processes that contribute to their mental health issues (Beck, 2020). A cognitive model typically includes beliefs, thoughts, emotions, and other relevant elements. The associations among these components provide a clear representation of the patient's maladaptive cognitive process (See Appendix B for an example cognitive model). Constructing a patient's cognitive model is a crucial step in CBT (Beck, 2020; Beck and Haigh, 2014; Wang et al., 2024), as therapists must first summarize these maladaptive structures before working to correct them. We believe that the process of cognitive model construction can be effectively enhanced with the assistance of LLMs, which could help in automatically identifying key components—such as beliefs and thought patterns—and their associations.

Existing work mostly focuses on cognitive distortion classification (Sharma et al., 2023c; Shreevastava and Foltz, 2021b; Chen et al., 2023b)–classifying the maladaptive thinking patterns in a patient's speech, which can be viewed as one component of cognitive modeling. However, existing datasets are typically limited by poor quality (Shreevastava and Foltz, 2021b; Chen et al., 2023b) or simplicity of the input speech. In our Level II tasks, we first propose a new dataset, CBT-CD, for cognitive distortion classification with en-

---

| Datasets | Inputs | Labels |
|---|---|---|
| CBT-CD | **Situation**: Our wedding was put off because his parents asked him to build a house for them 2 months before our wedding! They had a perfectly good house at the time they just wanted their dream house. **Thoughts**: I am a victim with no power in this situation. I must accept this behavior. I am too scared to leave this situation. I am not worthy of better. His parents hate me. His parents do not want us to get married. He may not want to marry me either. He loves his parents more than me. I will always be second in his life. His parents had no need for a house, and I know this for sure. I am aware of every aspect of this situation. | all-or-nothing thinking; personalization; mind reading |
| CBT-PC | **Situation**: I had an amazing childhood. When I was twelve in 2004, my father had to go to Iraq. My mother thought it would be best if she moved my brother and I back to the U.S., where we would have family support. I was very depressed because my dad was my hero and I blamed my mom for everything that went wrong. I felt like no one understood me... When my dad came back, he wanted a divorce from my mother. **Thoughts**: Everything was great until my mom messed everything up. Because of her, my brother and I had to leave our dad. We ended up living somewhere where no one liked me. It's her fault that I never felt like I fit in. Even when my dad came back, he didn't want us either - he wanted a divorce. I guess no one will ever want me in their life. I'll probably be alone forever. | helpless; unlovable |
| CBT-FC | **Situation**: My daughter was recently diagnosed as bipolar. If I say anything about seeking treatment, my daughter accuses me of not understanding her and what is happening to her. She is very paranoid and worries about her safety all the time. I need to know how to talk to her and what to do to get her into treatment. **Thoughts**: I am a bad mother. This is my fault. It is so shameful that my daughter has bipolar. If my daughter gets worse, then it will be my fault for not getting her into treatment. I need to do something. This is my responsibility. | I am incompetent; I am helpless; I am powerless, weak, vulnerable; I am bad - dangerous, toxic, evil |

Table 2: Examples from CBT-CD (cognitive distortion classification), CBT-PC (primary core belief classification), and CBT-FC (fine-grained core belief classification). See more examples in Appendix B.

|  | CBT-CD | CBT-PC | CBT-FC |
|---|---|---|---|
| # of examples | 146 | 184 | 112 |
| # of labels | 10 | 3 | 19 |
| Average situation length | 232.9 | 240.7 | 233.4 |
| Average thought length | 258.8 | 256.9 | 248.4 |
| Average ground truth labels | 2.5 | 1.9 | 3.8 |

Table 3: Statistics of three level II tasks.

hanced quality. We use the label set in (Shreevastava and Foltz, 2021b) of ten distortion types.

We then propose two new tasks centered on core beliefs. Core beliefs are at the center of a patient's cognitive model—deeply ingrained perceptions about oneself, others, and the world (Beck, 2020). Eliciting core beliefs is one of the most challenging stages of cognitive modeling, while the essential focus in later sessions is to correct such unhealthy, biased core beliefs. According to Beck (2020), there are three primary core beliefs: *helpless, unlovable,* and *worthless.* Each of these categories encompasses more specific, fine-grained core beliefs. For example, *I am incompetent* and *I am a loser* fall under *helpless.* In total, these fine-grained core beliefs are classified into 19 categories. In our Level II tasks, we propose new datasets CBT-PC for primary core belief classification, and CBT-FC for fine-grained core belief classification. Check Appendix A for detailed label definitions and distribution for all three datasets.

We use the patient speech in the TherapistQA dataset (Shreevastava and Foltz, 2021b), sourced from mental health-related posts in online forums, as our data source to construct CBT-CD, CBT-PC, and CBT-FC. Our annotation process is as follows: (1) From each original post, we instruct the expert

annotators to extract a segment that describes a situation along with the associated thoughts. Annotators may also create plausible imagined situations and thoughts to supplement the original segments if they are insufficiently detailed. (2) For cognitive distortions, we instruct the experts to annotate up to three key ones. For primary core beliefs and fine-grained core beliefs, we instruct the experts to annotate up to three primary beliefs. Under each identified primary core belief, we annotate up to three fine-grained beliefs.

We employed a total of six CBT experts from UpWork to conduct the annotation. Check Appendix H for our annotation interface. We conducted another round of cross-verification to filter out the low-quality ones. We end up with 146 examples for CBT-CD, 184 examples for CBT-PC, and 112 examples for CBT-FC. We also have another 20 examples for each task serving as other potential uses. Table 2 and 3 presents some statistics and examples of the three datasets. We employ two new experts to perform another round of annotation for each dataset to establish human expert performance references. We ended up with a weighted F-1 score of 49.1%, 77.7%, and 54.6% for CBT-CD, CBT-PC, and CBT-FC, respectively. Despite implementing a cross-verification step to filter out low-quality data, the results suggest that these tasks inherently involve a degree of subjectivity.

## 3.3 Level III: Therapeutic Response Generation

At this level, we take a further step to explore a longstanding and challenging question: ***Do LLMs have the ability to conduct effective therapeutic conversations?*** This ultimate capability assesses

| Category | Difficulty Level | Patient Speech | Reference Response |
|---|---|---|---|
| Negotiating a session agenda | Beginner | [Nervous] I don't think I'm ready for working on this today. | OK. We can revisit and possibly modify our plan for today. How about we first take a step back and explore your thinking about this? What thoughts are you noticing as we discuss the agenda? |
| Negotiating a session agenda | Intermediate | [Agitated] Wow. You won't believe what happened this week. It's a really long story . . . | It sounds like there's a lot on your mind, and I'd like to hear about it. Would it be OK to take a second to discuss our agenda for the day first, including where discussing this past week might fit in, as well as anything else you want to take up here today? I want to make sure that we budget our time accordingly. Shall we start with a homework check-in and then tackle the story? |
| Responding to therapeutic alliance ruptures | Advanced | [Anxious] I did the measure you asked me to fill out. Honestly, I think I might have rated you lower than usual, but I'm not sure I want to talk about it. | I was just noticing that your trust in me has gone down some. I wonder if you could help me appreciate what that's like for you? I'd far prefer that to just persisting with our agenda when you may have diminishing faith in it or me. |

Table 4: Three example deliberate practice exercises from CBT-DP. The patient speech in each exercise starts with an emotion label indicating the patient's emotional status.

CBT competence as a whole, beyond the specific tasks in levels I and II. In real practice, we must proceed with caution, using such applications to assist therapists rather than interacting directly with patients without supervision. Nonetheless, first and foremost, it is crucial to understand the LLMs' capability in providing effective responses.

Collecting a sufficient amount of real CBT session data is difficult due to privacy constraints. Taking the suggestions from professors specialized in CBT education (our co-authors), we propose using *Deliberate Practice* (Boswell and Constantino, 2022) as our level III task, CBT-DP. This methodology is traditionally used to assess graduate students' proficiency in CBT. Unlike simulated scenarios such as role-playing, deliberate practice involves a structured set of exercises focusing on responding to typical patient dialogue in CBT sessions. These exercises are divided into ten categories, such as *Working With Cognitions* and *Responding to Client Resistance*, addressing key aspects and challenges of conducting CBT. This approach ensures comprehensive coverage of critical and difficult situations likely to arise in real therapy. Generating high-quality responses to patient speech in CBT-DP can serve as an effective proxy for evaluating real-world CBT session efficacy.

Each category of CBT-DP contains approximately 15 exercises across three difficulty levels: beginner, intermediate, and advanced. We use the 156 exercises outlined in Boswell and Constantino (2022) as CBT-DP, where the authors provide one reference response for each exercise. Table 4 shows some example exercises with reference responses. Check Appendix A for all categories and statistics of all exercises. For model evaluations, we frame CBT-DP as a generation task. The patient's speech in each exercise serves as the input, and the model is expected to generate an appropriate response. We work with our co-authors (professors in social work) to propose a set of evaluation criteria, four fine-grained criteria for each category (Appendix A). As suggested by our co-authors, while Boswell and Constantino (2022) provide a set of reference responses, these should not be treated as the sole 'gold standard.' Consequently, matching-based automatic evaluation metrics may not provide an accurate measure for this task. Therefore, we employ experts to conduct pairwise comparisons between the model generation and reference response under each criterion.

## 4 Level I and II Experiments

### 4.1 Settings

We evaluated six popular LLMs' performance on basic CBT knowledge acquisition (level I) with CBT-QA and cognitive model understanding (level II) with CBT-CD, CBT-PC, and CBT-FC. The models include a close-sourced model GPT-4o[5] and five open-sourced models of sizes ranging from 7B to 405B: Mistral-v0.3-7B (Jiang et al., 2023), Gemma-2-9B (Gemma Team, 2024), Llama-3.1-8B, Llama-3.1-70B, and Llama-3.1-405B (Llama Team, 2024)[6]. We refer to models with sizes under 10B as small models while others as large models (including GPT-4o).

We cast the classification tasks of level II into multiple-choice questions, where each question can have more than one option. Prompt examples for each task of level I and II can be found in Appendix C. The inferencing temperature is set as 0.0 to eliminate the randomness of the model generations for reproducibility and data type bfloat16 is used for the model weights and activations.

---

[5]https://openai.com/index/gpt-4o-system-card/

[6]We always use the instruct version of the open-sourced models for all the experiments in this paper.

## 4.2 Results

Table 5 shows the performance of the LLMs. We report accuracy for CBT-QA and weighted precision, recall, and F1 score for the other three datasets. Here are the findings:

**Large LLMs are better at answering CBT knowledge questions.** Large models could achieve higher accuracies on CBT-QA than small models. This could be attributed to larger models storing more knowledge. However, Gemma-2-9B has a similar model size as Llama-3.1-8B and Mistral-v0.3-7B and it surpasses these two by a significant margin, even comparable to Llama3.1-70B. This may be caused by the distribution difference of the data used for training.

**Simply making the models larger could not enhance their understanding ability of the cognitive model.** Mistral-v0.3-7B and Llama-3.1-7B these two small models get the best F1 scores on CBT-CD and CBT-PC, outperforming models of large sizes. This gives us insight that when increasing the model size to pursue better general capabilities like reasoning, the expertise of mental health care should also receive attention.

**Current LLMs struggle with detecting fine-grained cognitive disorders or core beliefs.** CBT-CD and CBT-FC are very challenging for current LLMs. The models generally perform poorly on these two datasets which are very difficult even for professional therapists. Enhancing these abilities of LLMs could make the therapy process more efficient and accurate, and our datasets serve as valuable resources to propel potential advances.

To delve further, we examine the accuracies for questions of different knowledge types in CBT-QA and detailed F1 scores on labels in CBT-CD, CBT-PC, and CBT-FC. The results for two fine-grained classification tasks, CBT-CD and CBT-FC, are shown in Figure 1[7]. We could see that for CBT-QA, questions of *basic CBT knowledge and concepts* and *others* are easier to answer than questions of *case studies* and *practical CBT knowledge*. This could be the reason that questions of *case studies* and *practical CBT knowledge* require more understanding of how to use the corresponding knowledge, not just reciting them. For CBT-CD, the models all struggle to identify *mental filter*, *manification*, and *overgeneralization* disorders. And large models are not consistently outperforming small

models. For *all-or-nothing thinking* disorder, two small models have the best F1 scores. And for *emotional reasoning* disorder, GPT-4o and Llama3.1-70B have the best performance while Llama3.1-405B even falls behind Llama3.1-7B. For CBT-PC, small models are better for *helpless* core belief, and large models take the lead for *worthless*. For *unlovable* belief, their performance is not much different. For CBT-FC, models are struggling with detecting *I am needy* core belief.

We go through the prediction of Llama-3.1-405B on the two fine-grained classification datasets and present some typical ones in Appendix G. The main reasons for a wrong prediction include: 1) the model is not sensitive to the indicators of cognitive disorders or core beliefs, such as "People do think I'm..." of *overgeneralization* disorder or "There is nothing enjoyable in my life" of *I don't deserve to live* core belief; 2) The model focuses on several disorders or beliefs while neglecting the others; and 3) The model makes a decision without sufficient supporting conditions. For example, determining the *personalization* disorder can be divided into one attributing the negative behavior to some person and this person should be himself; one of the cases in CBT-FC only meets the first but is judged as *personalization* by the model.

## 5 Level III Experiments

### 5.1 Settings

We evaluate three LLMs' performance for CBT-DP (level III): Llama-3.1-8B, Llama-3.1-405B, and GPT-4o. For all generations, we keep the generation temperature at 0.7. For evaluation, as suggested by domain experts (§3.3), we conduct pairwise comparisons between model generations and the reference responses in (Boswell and Constantino, 2022) under our proposed criteria set (§3.3). To ensure a fair comparison, we prompt the LLMs to generate responses within similar lengths as the reference (see Appendix C).

### 5.2 Expert Evaluation

To minimize bias in the annotation process, we mixed and randomly shuffled the comparison pairs across all models. For each annotation pair, we randomized the order in which the model-generated and reference responses appeared in the annotation interface. This approach ensures that annotators remain unaware of the source of each response. We instruct the experts to label each pair with one

---

[7]The results for CBT-QA and CBT-PC are shown in Appendix D for space limit.

| Model | CBT-QA | CBT-CD | | | CBT-PC | | | CBT-FC | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | F1 | Accuracy | Recall | F1 | Accuracy | Recall | F1 |
| Mistral-v0.3-7B | 75.9 | 43.8 | 56.4 | **45.3** | 78.8 | 75.0 | 72.5 | 43.8 | 59.7 | 47.9 |
| Gemma-2-9B | 90.0 | 44.5 | 47.2 | 40.1 | 79.5 | 82.9 | 80.4 | 45.4 | 63.3 | 51.5 |
| Llama-3.1-8B | 78.2 | 39.7 | 59.9 | 44.1 | 75.2 | 93.9 | **82.5** | 37.0 | 62.8 | 44.8 |
| Llama-3.1-70B | 92.7 | 50.9 | 51.8 | 43.9 | 82.6 | 72.7 | 75.8 | 62.0 | 56.6 | 55.4 |
| Llama-3.1-405B | **95.0** | 49.4 | 44.2 | 43.3 | 85.3 | 70.1 | 75.3 | 53.1 | 68.7 | **58.2** |
| GPT-4o | 94.1 | 55.8 | 52.0 | 43.9 | 80.2 | 77.6 | 78.4 | 54.5 | 62.1 | 56.7 |
| Human | 90.7 | 51.1 | 48.1 | 49.1 | 76.0 | 79.4 | 77.6 | 53.0 | 57.3 | 54.6 |

Table 5: Performance of LLMs on basic CBT knowledge acquisition (CBT-QA) and cognitive model understanding (CBT-CD, CBT-PC, and CBT-FC). The precision, recall, and F1 are averaged by class portion.
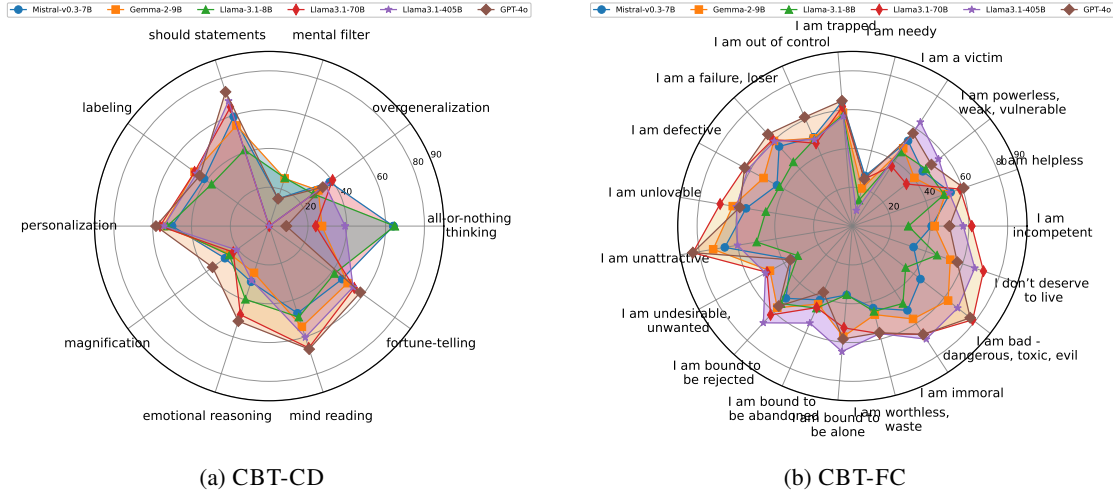


(a) CBT-CD

(b) CBT-FC

Figure 1: Detailed F1 scores of each label for CBT-CD and CBT-FC.

| Exercise | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Llama-3.1-405B | 0.07 | 0.06 | 0.21 | -0.24 | -0.19 | 0.00 | 0.18 | -0.31 | 0.07 | 0.00 | -0.01 |
| Llama-3.1-8B | -0.21 | -0.31 | -0.47 | -0.29 | -0.13 | 0.00 | -0.35 | 0.00 | -0.13 | -0.33 | -0.22 |
| GPT-4o | -0.50 | -0.50 | -0.13 | -0.24 | -0.31 | -0.53 | -0.06 | -0.44 | -0.13 | -0.40 | -0.32 |

Table 6: Overall rating: pairwise comparison of different models vs. reference across ten exercises.

of the following five labels: 1) A is much better than B, 2) A is slightly better than B, 3) A and B are about the same, 4) B is slightly better than A, and 5) B is much better than A. Labeling is done for each criterion defined, as well as an overall preference. We map the results to a scale from -2 to 2, where +2 indicates a strong preference for model-generated response.

The full annotation consists of 468 pairwise comparisons across all models, with each instance requiring approximately 3 minutes for the expert to finish. Given the expert's rate of $60 per hour, the annotation process takes around 24 hours to complete, resulting in a significant cost. We first have our expert co-author annotate a sample and write explanations of scoring. Then, we hired another

expert to finish the annotation of all data following the guidance of our co-author (Appendix H for our annotation interface).

### 5.3 Overall Result and Analysis

Here we first report the overall pairwise comparison results between each LLM vs. reference on 10 exercise categories in Table 6. Notably, the average results show that all three LLMs still lag behind humans by a significant margin, with Llama-3.1-405B achieving the best result while GPT-4o is the worst. Moreover, Llama-3.1-8B and GPT-4o never surpass human-written references in every single class of exercise, while Llama-3.1-405B shows some superior results in half of the exercise categories. Further insights reveal that Llama-3.1-405B is better because it *keeps engaging the client and ask-*

| Statement | [Frustrated] I know I said I would go to the event and start a conversation with new people, but I only managed to say "hello" to one person. |
|---|---|
| Llama-3.1-405B | It sounds like you're feeling disappointed with yourself for not meeting your own expectations. Can we take a closer look at what held you back from starting more conversations? Was it anxiety, fear of rejection, or something else? Let's break it down and identify some smaller, achievable steps for next time. |
| Human Reference | It seems like you might be a bit disappointed by this, but I see this as significant progress. It was a big accomplishment to show up to the event in the first place and then say hello to someone you don't know. Let's talk about what that was like. |

Table 7: Beginner exercise example: comparison between Llama-3.1-405B and human reference responses. The former is more structured, while the latter provides more emotional support.

*ing pertinent questions to guide them*, but GPT-4o tends to *correct the client's statement and keeps using professional CBT terminology to challenge the client*. And Llama-3.1-8B lies in the middle. Although this might be improved by adding additional limits in the prompt, we believe the current prompt starting with "you are a professional therapist using cognitive-behavioral therapy (CBT) in a session with a client." should already inform the model about its role. The complete prompt can be found in Table 22. We also draw the win-tie-loss figure in Appendix F.

Next, we discuss the overall score categorized by difficulty level in Figure 2. Surprisingly, we find that the models generally perform better when the difficulty level increases. Using Llama-3.1-405B as an example, we observe that the model's responses tend to be structured and rigid, often following a pattern such as repeating the client's statement, questioning reasons, and proposing possibilities. In contrast, human reference responses are more effective in affirming the patient's feelings and thoughts, and providing guidance *from the patient's perspective with a high degree of empathy and flexibility*, thereby establishing trust and rapport. At the beginner level, tasks tend to be logically simpler and more empathy-focused; consequently, model-generated answers are rated lower, as shown in Table 7.
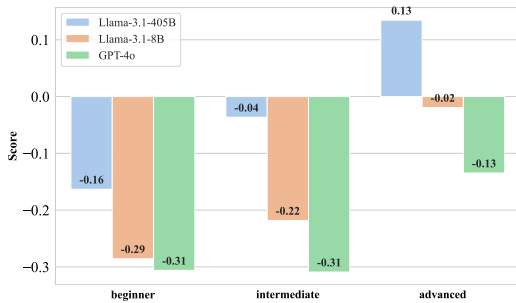


Figure 2: The overall pairwise comparison of different models vs. reference across difficulty level.

In addition, Table 8 shows the models' scores

| #Exe. | Metric | Model Results | | |
|---|---|---|---|---|
| | | Llama-3.1-405B | Llama-3.1-8B | GPT-4o |
| 2 | C1 | 0.34 | -0.25 | 0.25 |
| | C2 | 0.06 | 0.13 | -0.25 |
| | C3 | -0.38 | -0.13 | -0.63 |
| | C4 | 0.56 | 0.13 | 0.81 |

Table 8: Results breakdown of 4 criteria on questions from exercise 2, reported with the average score.

across 4 different criteria (C1-4) for exercises 2 (full results for all exercises in Appendix E). Apart from the overall rating, these criteria give fine-grained insights into the model's performance. We saw mixed results from the evaluation of exercise 2 (Establishing Goals). For example, in the comparison in appendix Table 24, criteria 1 is "Suggest CBT-consistent goals and tasks that align with an individualized CBT case formulation" and criteria 4 is "Emphasize concrete, actionable, and measurable goals". The human reference is much better in Criteria 1 and slightly worse in Criteria 4. Indeed, the model provides *concrete and actionable suggestions (aligns with criteria 4) while the reference focuses on CBT case formulation (aligns with criteria 1)*. However, for exercise 10 (Responding to Client Resistance), all criteria witness negative results on all models (see Appendix E), consistent with the overall rating. The reason can be summarized as the models tend to *dismiss the clients' experience and challenge their resistance rather than being responsive to patient needs as human therapists*. We show two such examples in appendix Table 25 and 26, clearly showing that the model tends to challenge the client while humans always respect clients' autonomy.

We provide a typical example comparing model responses in Table 29. In this case, the client perceives the therapist's suggestion of 'little rewards' as unreasonable. GPT-4o responds by directly using professional CBT terminology to correct the client's perspective, which might make the client feel uncomfortable. Llama-3.1-8B also addresses

the client's perspective but does so in a more gentle tone, discussing the underlying reasons for the issue with the client. Llama-3.1-405B delivers the best response. It uses an analogy to help the client view 'little rewards' from a different perspective and suggests collaboratively exploring the reasons behind the client's discomfort. This approach reduces the client's resistance and fosters greater willingness to engage in self-exploration and self-change.

## 6 Conclusion

In this paper, we propose a new benchmark, CBT-BENCH, to comprehensively evaluate the efficacy of assisting CBT with LLMs. CBT-BENCH contains three levels of tasks ranging from basic CBT knowledge acquisition to therapeutic response generation, totaling five new datasets constructed by domain experts. We experiment with representative LLMs and the results show that models achieve comparable or even better cognitive understanding ability, especially in QA tasks, but still lag behind humans in most generation tasks. We hope our new benchmark provides a meaningful testbed for future research in this area.

## 7 Limitations

In this work, we propose a set of novel tasks aimed at addressing key stages of CBT where LLMs could potentially assist professionals. We collaborated closely with domain experts to select tasks that are both important and challenging, and that are well-suited for AI enhancement. However, annotating datasets in highly specialized fields like mental health is very costly, which has limited the size of our datasets and constrained our ability to scale to additional tasks. There may be other stages in CBT that could also benefit from AI assistance, which we leave for future research. Additionally, for the level III task, we used deliberate practice as an effective proxy to evaluate performance, given the difficulty of acquiring real CBT session data due to privacy concerns. Bridging this gap by approximating real CBT sessions, while respecting ethical and privacy constraints, remains an important direction for future work.

## 8 Ethics Statement

**IRB (Institutional Review Board) Approval.** This project is approved by our Institutional Review Board (IRB). For all the data annotation and

evaluation, we work with our co-authors (clinical psychologists, professors, and social workers) and hire experts from UpWork[8]. Experts are defined as those with a graduate degree in clinical psychology, social work, or other related majors and have received at least 5 hours of CBT training. We hired a total of 8 US-based experts, each at an hourly rate of $60. For some experts, we agreed on a fixed payment based on an estimated completion time, calculated using the same hourly rate.

**Informed Consent.** All experts hired were 18 or older and provided informed consent (Check Appendix H for the consent forms). We did not assess any clinical outcomes. All data collected from the experts were de-identified and consented to be released for research purposes.

**System and Data Usages.** In this work, our primary goal is to systematically evaluate the capabilities of large language models (LLMs) in assisting Cognitive Behavioral Therapy (CBT) in order to provide insights for the future development of AI-based tools to support psychotherapy professionals. All data and systems developed in this work are intended exclusively for academic research purposes. They should not be used as real applications. We recognize the importance of human expertise and judgment in therapeutic settings and affirm that AI should be used to augment, not substitute, the essential role of trained professionals. Future developments of AI-based systems should not be used directly to interact with real patients without the supervision of professionals. We strongly advocate for rigorous ethical oversight in the deployment of AI-based systems in sensitive areas like mental health (see Appendix J for more discussion about ethics considerations). Furthermore, all data used in this study were either publicly available or collected with full transparency and consent, ensuring the privacy and confidentiality of individuals involved. All data collected from the experts were de-identified and consented to be released for research purposes.

## 9 Acknowledgment

We would like to thank Deborah J. Moon and Kate Hardy for providing CBT exam questions and useful discussions.

---

# References

Achini Adikari, Daswin De Silva, Harsha Moraliyage, Damminda Alahakoon, Jiahui Wong, Mathew Gancarz, Suja Chackochan, Bomi Park, Rachel Heo, and Yvonne Leung. 2022. Empathic conversational agents for real-time monitoring and co-facilitation of patient-centered healthcare. *Future Generation Computer Systems*, 126:318–329.

Aaron T Beck and Emily AP Haigh. 2014. Advances in cognitive theory and therapy: The generic cognitive model. *Annual review of clinical psychology*, 10(1):1–24.

Judith S Beck. 2020. *Cognitive behavior therapy: Basics and beyond*. Guilford Publications.

James F Boswell and Michael J Constantino. 2022. *Deliberate practice in cognitive behavioral therapy*. American Psychological Association.

Siyuan Chen, Mengyue Wu, Kenny Q Zhu, Kunyao Lan, Zhiling Zhang, and Lyuchun Cui. 2023a. LLM-empowered chatbots for psychiatrist and patient simulation: Application and evaluation. *arXiv [cs.CL]*.

Zhiyu Chen, Yujie Lu, and William Wang. 2023b. Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4295–4304.

Shuang Gao, Vince D Calhoun, and Jing Sui. 2018. Machine learning in major depression: From classification to treatment outcome prediction. *CNS neuroscience & therapeutics*, 24(11):1037–1052.

Gemma Team. 2024. Gemma 2: Improving open language models at a practical size. *arXiv [cs.CL]*.

Declan Grabb, Max Lamparth, and Nina Vasan. 2024. Risks from language models for automated mental healthcare: Ethics and structure for implementation. *arXiv [cs.CY]*.

Shang-Ling Hsu, Raj Sanjay Shah, Prathik Senthil, Zahra Ashktorab, Casey Dugan, Werner Geyer, and Diyi Yang. 2023. Helping the helper: Supporting peer counselors via ai-empowered practice and feedback. *arXiv preprint arXiv:2305.08982*.

Shaoxiong Ji, Tianlin Zhang, Kailai Yang, Sophia Ananiadou, and Erik Cambria. 2023. Rethinking large language models in mental health applications. *arXiv [cs.CL]*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. *arXiv [cs.CL]*.

Dae-Young Kim, Rebecca Hwa, and Muhammad Mahbubur Rahman. 2024. MhGPT: A lightweight generative pre-trained transformer for mental health text analysis. *arXiv [cs.CL]*.

Willem Kuyken, Christine A Padesky, and Robert Dudley. 2011. *Collaborative case conceptualization: Working effectively with clients in cognitive-behavioral therapy*. Guilford Press.

June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. 2023. ChatCounselor: A large language models for mental health support. *arXiv [cs.CL]*.

Llama Team. 2024. The llama 3 herd of models. *arXiv [cs.AI]*.

Ryan Louie, Ananjan Nandi, William Fang, Cheng Chang, Emma Brunskill, and Diyi Yang. 2024. Roleplay-doh: Enabling domain-experts to create LLM-simulated patients via eliciting and adhering to principles.

Mounica Maddela, Megan Ung, Jing Xu, Andrea Madotto, Heather Foran, and Y-Lan Boureau. 2023. Training models to generate, recognize, and reframe unhelpful thoughts. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Stroudsburg, PA, USA. Association for Computational Linguistics.

Saskia Senn, ML Tlachac, Ricardo Flores, and Elke Rundensteiner. 2022. Ensembles of bert for depression classification. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 4691–4694. IEEE.

Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2021. Towards facilitating empathic conversations in online mental health support: A reinforcement learning approach. In *Proceedings of the Web Conference 2021*, pages 194–205.

Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2023a. Human–ai collaboration enables more empathic conversations in text-based peer-to-peer mental health support. *Nature Machine Intelligence*, 5(1):46–57.

Ashish Sharma, Kevin Rushton, Inna Lin, David Wadden, Khendra Lucas, Adam Miner, Theresa Nguyen, and Tim Althoff. 2023b. Cognitive reframing of negative thoughts through human-language model interaction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Stroudsburg, PA, USA. Association for Computational Linguistics.

Ashish Sharma, Kevin Rushton, Inna Lin, David Wadden, Khendra Lucas, Adam Miner, Theresa Nguyen, and Tim Althoff. 2023c. Cognitive reframing of negative thoughts through human-language model interaction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9977–10000.

Ashish Sharma, Kevin Rushton, Inna Wanyin Lin, Theresa Nguyen, and Tim Althoff. 2024. Facilitating self-guided mental health interventions through human-language model interaction: A case study of cognitive restructuring. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, volume 21, pages 1–29, New York, NY, USA. ACM.

Sagarika Shreevastava and Peter Foltz. 2021a. Detecting cognitive distortions from patient-therapist interactions. In *Proceedings of the Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access*, Stroudsburg, PA, USA. Association for Computational Linguistics.

Sagarika Shreevastava and Peter Foltz. 2021b. Detecting cognitive distortions from patient-therapist interactions. In *Proceedings of the Seventh Workshop on Computational Linguistics and Clinical Psychology: Improving Access*, pages 151–158.

Bichen Wang, Pengfei Deng, Yanyan Zhao, and Bing Qin. 2023. C2D2 dataset: A resource for the cognitive distortion analysis and its impact on mental health. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10149–10160, Stroudsburg, PA, USA. Association for Computational Linguistics.

Ruiyi Wang, Stephanie Milani, Jamie C Chiu, Shaun M Eack, Travis Labrum, Samuel M Murphy, Nev Jones, Kate Hardy, Hong Shen, Fei Fang, et al. 2024. Patient-{\Psi}: Using large language models to simulate patients for training mental health professionals. *arXiv preprint arXiv:2405.19660*.

Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyan Kuang, and Sophia Ananiadou. 2023. Towards interpretable mental health analysis with large language models. *arXiv [cs.CL]*.

Kailai Yang, Tianlin Zhang, Ziyan Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. MentaLLaMA: Interpretable mental health analysis on social media with large language models. In *Proceedings of the ACM on Web Conference 2024*, volume 35 of *WWW '24*, page 4489–4500. ACM.

# A Label Distribution and Definitions

| Belief | #Label | Definition |
|---|---|---|
| all-or-nothing thinking | 65 | Also called black-and-white, polarized, or dichotomous thinking. You view a situation in only two categories instead of on a continuum. |
| mind reading | 47 | You believe you know what others are thinking, failing to consider other, more likely possibilities. |
| fortune-telling | 44 | Also called fortune-telling. You predict the future negatively without considering other, more likely outcomes. |
| personalization | 42 | You believe others are behaving negatively because of you, without considering more plausible explanations for their behavior. |
| emotional reasoning | 36 | You think something must be true because you "feel" (actually believe) it so strongly, ignoring or discounting evidence to the contrary. |
| overgeneralization | 32 | You make a sweeping negative conclusion that goes far beyond the current situation. |
| labeling | 29 | You put a fixed, global label on yourself or others without considering that the evidence might more reasonably lead to a less extreme conclusion. |
| should statements | 28 | Also called imperatives. You have a precise, fixed idea of how you or others should behave, and you overestimate how bad it is that these expectations are not met. |
| magnification | 25 | When you evaluate yourself, another person, or a situation, you unreasonably magnify the negative and/or minimize the positive. |
| mental filter | 21 | Also called selective abstraction. You pay undue attention to one negative detail instead of seeing the whole picture. |

Table 9: Label distribution and definitions of CBT-CD.

| Primary Belief | #Label | Definition |
|---|---|---|
| helpless | 162 | being ineffective—in getting things done, self-protection, and/or measuring up to others |
| unlovable | 101 | having personal qualities resulting in an inability to get or maintain love and intimacy from others |
| worthless | 81 | being an immoral sinner or dangerous to others |

Table 10: Label distribution and definitions of CBT-PC.

Table 11: Label distribution of CBT-FC.

| Primary Core Belief | Fine-grained Core Belief | #Label |
|---|---|---|
| helpless | I am incompetent | 16 |
| | I am helpless | 36 |
| | I am powerless, weak, vulnerable | 38 |
| | I am a victim | 29 |
| | I am needy | 5 |
| | I am trapped | 28 |
| | I am out of control | 22 |
| | I am a failure, loser | 22 |
| | I am defective | 31 |
| unlovable | I am unlovable | 25 |
| | I am unattractive | 6 |
| | I am undesirable, unwanted | 27 |
| | I am bound to be rejected | 29 |
| | I am bound to be abandoned | 24 |
| | I am bound to be alone | 21 |
| worthless | I am worthless, waste | 31 |
| | I am immoral | 10 |
| | I am bad - dangerous, toxic, evil | 14 |
| | I don't deserve to live | 8 |

Table 12: Statistics of CBT-DP.

| Category | #Exercise | #Beginner | #Intermediate | #Advanced |
|---|---|---|---|---|
| Explaining the Treatment Rationale for Cognitive Behavioral Therapy | 14 | 4 | 6 | 4 |
| Establishing Goals | 16 | 5 | 6 | 5 |
| Negotiating a Session Agenda | 15 | 5 | 5 | 5 |
| Assigning and Reviewing Between-Session Activities | 15 | 5 | 7 | 5 |
| Working With Cognitions | 16 | 5 | 5 | 6 |
| Working With Behaviors | 15 | 5 | 5 | 5 |
| Working With Emotions | 17 | 5 | 5 | 7 |
| Adherence Flexibility | 16 | 5 | 6 | 5 |
| Responding to Therapeutic Alliance Ruptures | 15 | 5 | 5 | 5 |
| Responding to Client Resistance | 15 | 5 | 5 | 5 |
| All | 156 | 49 | 55 | 52 |

| Category | Criteria |
|---|---|
| Explaining the Treatment Rationale for Cognitive Behavioral Therapy | Criteria 1: Validate the client's experience,<br>Criteria 2: Explain the logic of how CBT can be used to address concerns,<br>Criteria 3: Instill hope for using CBT effectively,<br>Criteria 4: Set appropriate expectations for the nature and impact of CBT |
| Establishing Goals | Criteria 1: Suggest CBT-consistent goals and tasks that align with an individualized CBT case formulation,<br>Criteria 2: Invite the client's input on and agreement with the goals and tasks,<br>Criteria 3: Demonstrate flexibility,<br>Criteria 4: Emphasize concrete, actionable, and measurable goals |
| Negotiating a Session Agenda | Criteria 1: Suggest a CBT-consistent agenda that aligns with the CBT case formulation,<br>Criteria 2: Invite the client's input on and agreement with the agenda,<br>Criteria 3: Demonstrate flexibility,<br>Criteria 4: Maintain some degree of frame or structure (i.e., the agenda can shift, but generally maintain a more or less explicit attempt to delineate expectations for session goals and tasks) |
| Assigning and Reviewing Between-Session Activities | Criteria 1: Provide a basic rationale for the homework assignment,<br>Criteria 2: Tailor the homework assignment to the client's concerns and needs,<br>Criteria 3: Ensure that there is adequate agreement on and understanding of the homework assignment,<br>Criteria 4: Encourage and validate good-faith efforts |
| Working With Cognitions | Criteria 1: Encourage and facilitate client self-reflection on thoughts and beliefs,<br>Criteria 2: Maintain an open, exploratory stance,<br>Criteria 3: Orient the client's attention to thoughts or connections between thoughts and other experiences,<br>Criteria 4: Emphasize cognitive flexibility, rather than simple thought replacement, and refrain from implying that a particular thought is "right or wrong" |
| Working With Behaviors | Criteria 1: Maintain a collaborative and curious stance,<br>Criteria 2: Be clear when suggesting a plan of action and expectations,<br>Criteria 3: Be positive and supportive of both big and small positive changes,<br>Criteria 4: Appeal to relevant core learning concepts, including conditioning, reinforcement, the importance of environment, antecedents, and consequences |
| Working With Emotions | Criteria 1: Empathically inquire about clients' emotional experience,<br>Criteria 2: Actively listen and maintain a supportive tone to clients' emotional disclosures,<br>Criteria 3: Model tolerance of affect and an approach orientation to clients' strong emotional experience,<br>Criteria 4: Use psychoeducation to support the importance of clients experiencing versus avoiding their emotion and the maladaptive consequences of negative reinforcement |
| Adherence Flexibility | Criteria 1: Maintain a collaborative and transparent stance,<br>Criteria 2: Demonstrate empathy,<br>Criteria 3: Demonstrate openness to explore before determining if a modification is indeed warranted,<br>Criteria 4: Maintain consistency with a broad CBT orientation |
| Responding to Therapeutic Alliance Ruptures | Criteria 1: Momentarily step away from the CBT change agenda,<br>Criteria 2: Warmly invite the client to discuss their subjective experience,<br>Criteria 3: Empathize with the client's thoughts and feelings and invite further disclosure of unhelpful or invalidating occurrences,<br>Criteria 4: "Disarm" the client's negative affect toward you or CBT by finding some truth in their disclosures |
| Responding to Client Resistance | Criteria 1: Momentarily step away from the CBT change agenda,<br>Criteria 2: Empathically explore the client's experience of treatment and you,<br>Criteria 3: Validate clients' experience and "roll with" versus challenging their resistance,<br>Criteria 4: Support clients' autonomy and elicit their motivation for pursuing valued directions |

Table 13: Evaluation criteria of CBT-DP.

# B Examples

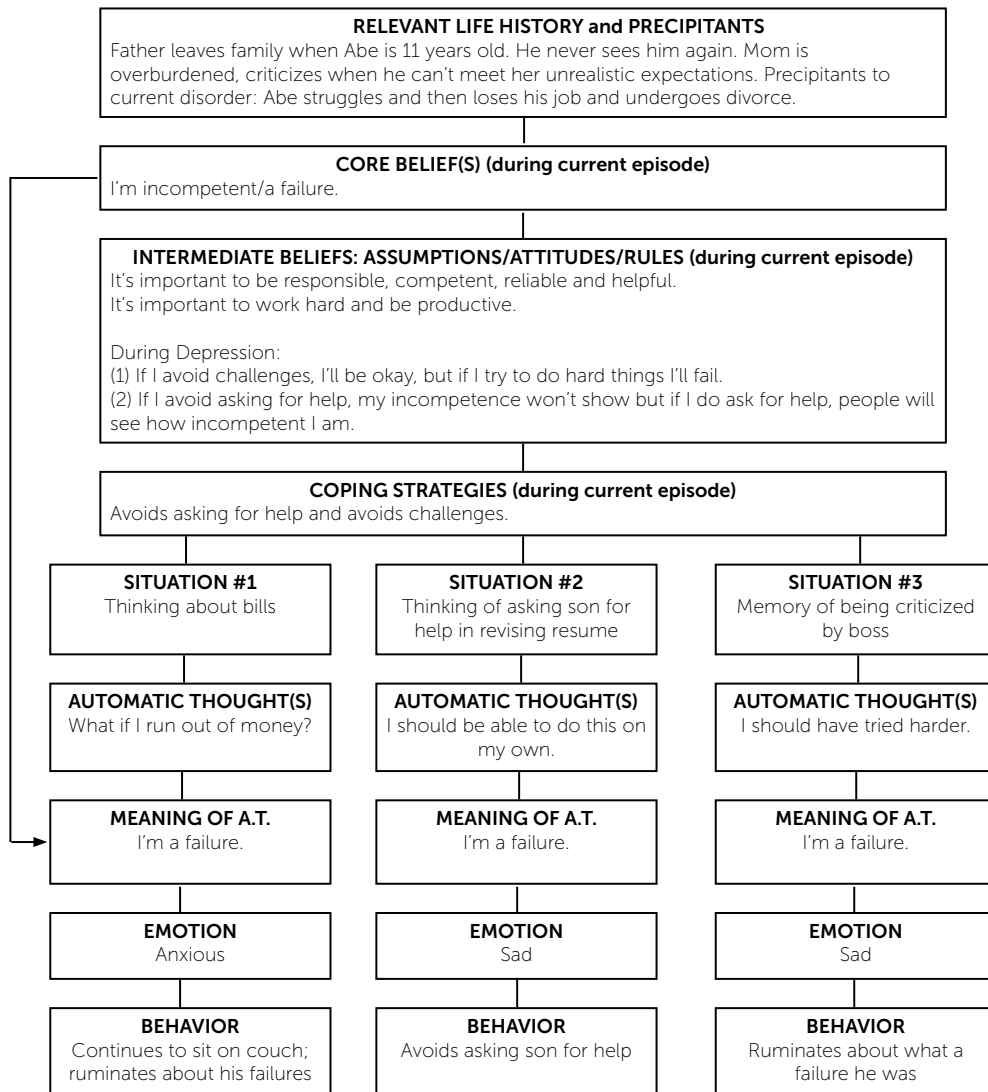| Knowledge Types | Example QA Pairs from CBT-QA |
|---|---|
| Basic CBT knowledge and concepts | Intermediate beliefs consist of rules, attitudes, and conditional assumptions. Which of the following fits the best as an attitude?<br>**A.** "Failing is bad", **B.** "I must not fail",<br>**C:** "If I try, I will fail", **D.** "I am a failure" |
| | What type of cognitive error involves shrinking something to make it seem less important?<br>A. Selective Abstraction, **B.** Mental filtering,<br>**C.** Overgeneralization, **D.** Self-defeated thinking |
| | Which strategy should you employ the most often when you want a client to come to a new understanding?<br>**A.** Guided discovery, **B.** Direct advice,<br>**C:** Indirect advice, **D.** Disclosing what strategy works the best for you |
| | In the assessment session, you should create a conceptualization/formulation of the client's problem with them.<br>When using CBT, which of the following is the most important for a conceptualization/formulation to include?<br>**A.** How past developmental events impact the problem, **B.** How thoughts and behaviors relate to the problem,<br>**C:** How current supports may limit the client's functioning, **D.** How genetic vulnerabilities result in current problems |
| | CBT helps clients identify links between:<br>**A.** Health, thoughts, **B.** Emotions, thoughts, decision-making,<br>**C:** Thoughts, emotions, behavior, **D.** Antecedent, behavior, consequence |
| Practical CBT knowledge | How might a social worker use the principles of CBT approach when working with a client who is experiencing suicidal thoughts?<br>**A.** By encouraging the client to suppress their thoughts and focus on positive affirmations,<br>**B.** By providing temporary shelter and financial assistance to alleviate immediate stressors,<br>**C.** By collaboratively identifying and challenging distorted thought patterns contributing to suicidal thoughts,<br>**D.** No options align with any practices considered to be CBT |
| | A primary concern during the final stages of a cognitive behavioral group is which of the following?<br>**A.** promoting transfer of learning, **B.** developing a therapeutic contract,<br>**C.** role-playing various situations **D.** establishing baseline data |
| | Early in treatment you advise clients that the way to get better is?<br>**A.** By making small changes in their thinking and behaviors each day, **B.** By making large changes to their core beliefs,<br>**C.** By gaining insight regarding their developmental history **D.** By changing their emotions |
| | When helping clients evaluate automatic thoughts, therapists should generally help clients evaluate which aspects of those thoughts?<br>**A.** accuracy and/or intensity, **B.** intensity and/or utility,<br>**C.** accuracy and/or utility |
| | Which of the following typically improve depressive symptoms (select all that apply)?<br>**A.** Interpersonal interactions, **B.** Activities that have the potential for pleasure,<br>**C.** Activities that have the potential for a sense of mastery, **D.** Repeatedly thinking about the reasons one is depressed |
| Case studies | The client says "I feel like he doesn't care about me". "I feel like he doesn't care about me" is a?<br>**A.** Thought, **B.** Emotion,<br>**C.** Physiological reaction/sensation, **D.** Feeling |
| | A manager tells you (the therapist) that the new employee, who has been doing an excellent job, requested a meeting with him.<br>The manager anticipates that the new employee plans to resign. Which is the best description of the manager's cognitive error?<br>**A.** Thinking from an all-or-nothing perspective, **B.** Always thinking the worst will occur without considering positive outcomes,<br>**C.** Viewing only selected negative evidence while editing out positive aspects, **D.** Undervaluing the positive significance of an event |
| | The client's mood has been improving since they've been attending treatment. They could attribute this improvement to many different things.<br>Which of the below do you want them to attribute it to the most?<br>**A.** The positive changes they have been making in their thinking and behaviors, **B.** The support their therapist has given them,<br>**C.** Events in their environment that have changed, **D.** The randomness of their mood |
| | Ms. T. has a number of long-standing negative attitudes such as "I'll never succeed" and "I have to be perfect to be accepted."<br>No matter how hard she tries she always seems to think that she is "not measuring up."<br>What cognitive therapy procedures do you think might help her?<br>**A.** Listing advantages and disadvantages, **B.** Breaking out,<br>**C.** Listing schemas, **D.** Thought recording |
| | A client says, "I don't feel like I'm doing a good job at work"<br>What question can you ask next to start the downward arrow technique to identify their core belief?<br>**A.** "How did your parents give feedback when you were younger?", **B.** "That's not true, you're a great worker!",<br>**C.** "What are you worried that might mean?", **D.** "What is the evidence for and against this thought?" |
| Others | Among the following principles, which one does not align with the core values of social work practice?<br>**A.** Respect for the inherent dignity and worth of every individual,<br>**B.** Promotion of client self-determination and autonomy,<br>**C.** Maintenance of strict neutrality and avoidance of taking sides,<br>**D.** Commitment to social justice and advocacy for marginalized populations |
| | Which of the following best describes the collaborative nature of a client and a therapist in social work practice?<br>**A.** The therapist provides solutions to the client's problems without considering the client's perspective,<br>**B.** The therapist and client work together as partners to identify goals and develop strategies for change,<br>**C.** The therapist imposes their own values and beliefs onto the client without discussion,<br>**D.** The therapist directs the client's treatment without input or collaboration from the client |
| | In the field of social work, what is the primary goal of intervention?<br>**A.** Providing financial assistance to individuals in need,<br>**B.** Offering temporary shelter to homeless populations,<br>**C.** Enhancing the well-being and social functioning of individuals, families, and communities,<br>**D.** Administering psychological assessments to diagnose mental health disorders |
| | A social worker needs to address the specific needs of oppressed populations, such as marginalized racial and ethnic groups.<br>What would be the best approach through a CBT framework?<br>**A.** Allowing others in the community to take part and take the lead. Mainly because as social workers we should allow others to take charge,<br>**B.** By encouraging clients to conform to dominant cultural norms to reduce societal backlash,<br>**C.** By validating that such issues and problems exist, and then seeing if your agency has any funding to help purchase food for families impacted,<br>**D.** By validating clients' experiences of discrimination and integrating culturally sensitive interventions |
| | How might therapists incorporate LGBTQ-affirmative approaches to address the unique challenges faced by LGBTQ individuals?<br>**A.** By avoiding discussions about clients' LGBTQ identities to focus solely on symptom management,<br>**B.** By validating the identities of LGBTQ-identified individuals through demonstrating acceptance and affirming their experiences,<br>**C.** By challenging clients' LGBTQ identity as a potential cause of their mental health struggles,<br>**D.** By encouraging clients to suppress their LGBTQ identity to avoid social stigma |

Table 14: Example QA pairs of CBT-QA.

| Inputs | Labels |
|---|---|
| **Situation**: I have been with my girlfriend for 5 years and we have a history of a strong sexual intimacy connection. This has been getting worse over time. **Thoughts**: It is my girlfriend's sex drive that is the problem. | overgeneralization; labeling |
| **Situation**: I am struggling at the moment and yet I am aware that I should be feeling good. **Thoughts**: What I am feeling and going through are unacceptable and wrong. There is a right and wrong way to go through this and I am doing it the wrong way. | should statements; emotional reasoning |
| **Situation**: The individual has recently discovered that their mom has cheated on their dad four times in the past. Each time, the parents attended counseling, where the mom initially denied her infidelity but eventually admitted to it, blaming their dad for not meeting her needs. **Thoughts**: I cannot trust anyone, why is everyone a liar. My mom has betrayed me. Of course something like this is happening to me, it only happens to me. | all-or-nothing thinking; overgeneralization; emotional reasoning |
| **Situation**: I experienced a recent setback with my OCD symptoms that had been in remission **Thoughts**: I have no longer won the battle against my OCD | all-or-nothing thinking; magnification |
| **Situation**: My anxiety used to be worse a couple of years ago, but now im just paranoid. Im always thinking that people think im a loser, and I wont leave the house without makeup, because i am horrified about what other people think about me. **Thoughts**: Because my anxiety felt worse before, this must be a normal level of anxiety and paranoia, so my conclusions make more sense. People do think I'm a loser or that I'm ugly. I at least wear makeup so they can't see how ugly I am, and if I avoid people altogether, then they won't know how much of a loser I am or how boring I am. I don't like who I am, and whatever people think of me is probably way worse than I think it is. | overgeneralization; emotional reasoning; mind reading |

Table 15: Examples from CBT-CD (Cognitive Distortion Classification).

**BECK**
INSTITUTE

# (TRADITIONAL) COGNITIVE CONCEPTUALIZATION DIAGRAM EXAMPLE

Name: _____ Date: _____ Diagnosis: _____

**RELEVANT LIFE HISTORY and PRECIPITANTS**
Father leaves family when Abe is 11 years old. He never sees him again. Mom is overburdened, criticizes when he can't meet her unrealistic expectations. Precipitants to current disorder: Abe struggles and then loses his job and undergoes divorce.

**CORE BELIEF(S) (during current episode)**
I'm incompetent/a failure.

**INTERMEDIATE BELIEFS: ASSUMPTIONS/ATTITUDES/RULES (during current episode)**
It's important to be responsible, competent, reliable and helpful.
It's important to work hard and be productive.

During Depression:
(1) If I avoid challenges, I'll be okay, but if I try to do hard things I'll fail.
(2) If I avoid asking for help, my incompetence won't show but if I do ask for help, people will see how incompetent I am.

**COPING STRATEGIES (during current episode)**
Avoids asking for help and avoids challenges.

| **SITUATION #1** Thinking about bills | **SITUATION #2** Thinking of asking son for help in revising resume | **SITUATION #3** Memory of being criticized by boss |
|---|---|---|
| **AUTOMATIC THOUGHT(S)** What if I run out of money? | **AUTOMATIC THOUGHT(S)** I should be able to do this on my own. | **AUTOMATIC THOUGHT(S)** I should have tried harder. |
| **MEANING OF A.T.** I'm a failure. | **MEANING OF A.T.** I'm a failure. | **MEANING OF A.T.** I'm a failure. |
| **EMOTION** Anxious | **EMOTION** Sad | **EMOTION** Sad |
| **BEHAVIOR** Continues to sit on couch; ruminates about his failures | **BEHAVIOR** Avoids asking son for help | **BEHAVIOR** Ruminates about what a failure he was |

Figure 3: An example cognitive model from Beck (2020).

| Inputs | Labels |
|---|---|
| **Situation**: I've been dating this guy for 2 years, and when things are good, they're great. We have deep conversations about various topics, share laughs, and seem to have a strong connection. However, every few months, he suddenly becomes distant. The usual daily communication slows down, his tone becomes cold, and there's a noticeable disconnect. After giving it some space and bringing it up gently, he admits to withdrawing and usually provides a reason for his behavior, such as feeling disrespected or abandoned. We address the issues, talk about feelings and needs, and find better ways to communicate. He has a therapist, and I encourage him to work on these issues, but the cycle repeats. **Thoughts**: I am not enough for him. Why do I always get hurt like this. I always get stuck in these endless cycles, and I am unlovable. | helpless; unlovable |
| **Situation**: I would move out, but I don't have the money for that nor do I have the job. Whenever I get I job, I take too many sick days out because of the bruises and scarring on my face are hard to hide. **Thoughts**: I can't move out until I have a better financial situation. I must remain in this situation until I have more money. The only way out is to be more financially secure. I cannot let other people know or help me. I have to sick days to hide the bruising on my face. Other people can't find out what's happening. I must solve this all on my own. | helpless; unlovable |
| **Situation**: Ever since, I have known in the back of my mind that he still does this, and have gotten up in the middle of the night a few times to find him masturbating and video chatting publicly or with females who are at least scantily clad, if not nude and masturbating themselves. When I confront him about it, he says he has a public sex fetish and needs to fulfill it, believing that it helps keep him „faithful‰ to me since I am no longer interested. **Thoughts**: He calls that faithful? It doesn't feel faithful. But I don't know if I even have the right to feel the way I do about it – after all, that's pretty much how our relationship started. Am I even allowed to be upset? I'm the one that changed, so it feels like it's my fault. I shouldn't be blaming him. He tried to get me involved, but when I said I wasn't into it, he didn't push, he just continued on doing what made him happy. Plus, he's not doing this to hurt me, so I shouldn't feel hurt by it. I should be happy for him that he is feeling fulfilled. But I feel sad, like he doesn't want me anymore. | helpless; unlovable; worthless |
| **Situation**: Do I have schizophrenia or something **Thoughts**: I am scared and don't know what is happening. I fear this could be something very serious. There is something wrong with me. I don't know what to do. | helpless |
| **Situation**: I was born with Attention Deficit Hyperactivity Disorder, O.C.D. and I am also transgendered. As a result, I was a very difficult child to raise. There was constant fighting between myself and my parents. At the age of ten, I was sent away to an institution for two years. While there, I was physically, emotionally, and sexually abused. I couldn't talk to anyone because I had no trust in adults. **Thoughts**: I was born defective and wrong, so that's why my parents had such a hard time raising me. I was hard to love because I was so different. It was my fault that they had to send me away; I needed to be fixed. Then, I was abused while I was away but I couldn't talk to anyone about it because I was scared that if I said anything, I'd be sent even further away. | helpless; unlovable; worthless |

Table 16: Examples from CBT-PC (Primary Core Belief Classification).

| Inputs | Labels |
|---|---|
| **Situation**: My relationship has always been in trouble because of my wife's suspicious nature. **Thoughts**: - This is all my wife's fault, not mine. If she stopped being suspicious, then all of our marital problems would be solved. There is nothing I can do about this; she is the one who needs to change. | I am incompetent; I am powerless, weak, vulnerable; I am a victim |
| **Situation**: My relationship with my mom has deteriorated due to a lie I told back in February. I took full responsibility for it, but I'm aware that my actions, such as hiding my depression from her and seeing a psychologist without her knowledge, were not acceptable. I felt that she wouldn't understand because she's always belittled my emotions, accusing me of exaggerating and seeking attention whenever I express myself emotionally. This led me to feel ashamed and afraid of facing her criticism, so I chose to keep my struggles hidden. **Thoughts**: I'm a failure for lying and hiding things from my mom. She'll never understand or accept me for who I am. I'm just seeking attention by expressing my feelings. | I am trapped; I am a failure, loser; I am unlovable; I am undesirable, unwanted; I am bound to be rejected; I am worthless, waste |
| **Situation**: ...although we all had a lovely day I noticed he was distant... **Thoughts**: I am confused by this situation, I do not know how to interpret this or respond. I cannot trust myself. I cannot trust my own thoughts about this. | I am incompetent; I am powerless, weak, vulnerable |
| **Situation**: I have been backwards... I was forced... **Thoughts**: My sleep habits are abnormal. I should feel conscientious and weird about my sleep patterns. I may be judged for being different. I regret or resent having been forced. | I am defective |
| **Situation**: I never told anyone the truth about it or should I say I never admitted it to anyone... I swore I would never cheat on my husband because I never forgave myself for cheating on him with the married guy... I hated myself for that. **Thoughts**: I feel guilt for what I did. I blame myself for what happened. I am a bad person. I regret what I did. I wanted to make amends for what I did. I am only worthy of love if I hide my past mistakes. I feel shame for what I did. | I am bound to be rejected; I am bound to be abandoned; I am immoral |

Table 17: Examples from CBT-FC (Fine-grained Core Belief Classification).

## C  Task Prompts

You are taking a CBT exam and doing multiple-choice questions. Each question has only one right choice.
Which core principle underlies Cognitive Behavioral Therapy?
a: Acceptance and Commitment
b: Mindfulness
c: The influence of desires
d: The interconnection between thoughts, feelings, and behaviors
Answer:
Please only output the letter corresponding to the choice.

Table 18: Prompt example for CBT Exam QA.

You are a CBT therapist and now need to determine the cognitive distortions of a patient from his current situation and thoughts.
Each patient may have **up to 3** distortions.
Situation: I'm depressed but nobody knows it. I do not have any friends. This started at age 11.
Thoughts: I cannot make friends and have no one to hang out with. Therefore, I am always going to be alone and depressed.
Question: what distortions this patient has?
Choices:
a: all-or-nothing thinking
b: overgeneralization
c: mental filter
d: should statements
e: labeling
f: personalization
g: magnification
h: emotional reasoning
i: mind reading
j: fortune-telling
Answer:
Please only output the letters corresponding to the choices. Multiple choices should be separated by a comma.

Table 19: Prompt example for Cognitive Distortion Classification.

You are a CBT therapist and now need to determine the major core beliefs of a patient from his current situation and thoughts.
Each patient may have multiple core beliefs.
Situation: I have a history of being hurt in relationships so decided to take a break from dating. I am now in a relationship but keep hurting him and taking everything out on him.
Thoughts: I am hurting my boyfriend because I have been hurt in the past. I have so many problems.
Question: what core beliefs this patient has?
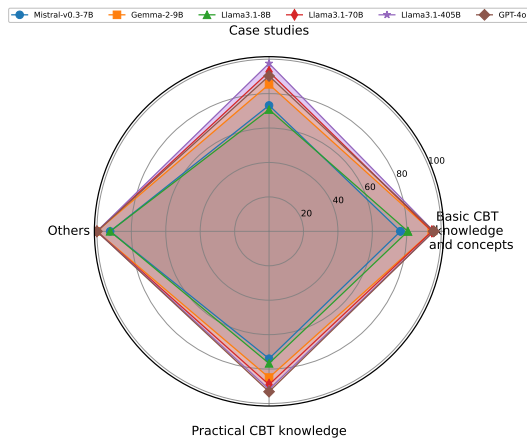a: helpless
b: unlovable
c: worthless
Answer:
Please only output the letters corresponding to the choices. Multiple choices should be separated by a comma.

Table 20: Prompt example for Primary Core Belief Classification.

You are a CBT therapist and now need to determine the fine-grained beliefs of a patient from his current situation and thoughts. Each patient may have **up to 9** fine-grained beliefs. Now answer the following question:

Situation: The individual has recently discovered that their mom has cheated on their dad four times in the past. Each time, the parents attended counseling, where the mom initially denied her infidelity but eventually admitted to it, blaming their dad for not meeting her needs.

Thoughts: I cannot trust anyone, why is everyone a liar. My mom has betrayed me. Of course something like this is happening to me, it only happens to me.

Question: what fine-grained beliefs this patient has?

Choices:

a: I am incompetent
b: I am helpless
c: I am powerless, weak, vulnerable
d: I am a victim
e: I am needy
f: I am trapped
g: I am out of control
h: I am a failure, loser
i: I am defective
j: I am unlovable
k: I am unattractive
l: I am undesirable, unwanted
m: I am bound to be rejected
n: I am bound to be abandoned
o: I am bound to be alone
p: I am worthless, waste
q: I am immoral
r: I am bad - dangerous, toxic, evil
s: I don't deserve to live

Answer:

Please only output the letters corresponding to the choices. Multiple choices should be separated by a comma.

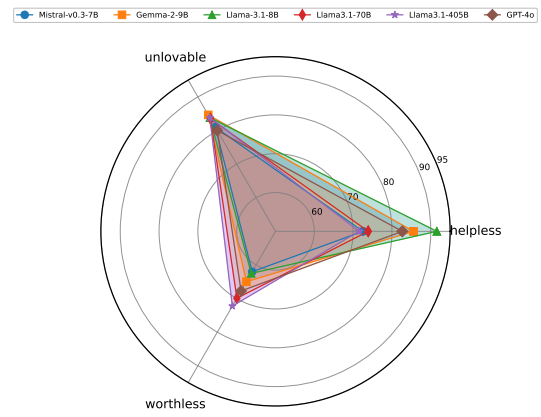Table 21: Prompt example for Fine-grained Core Belief Classification.

Now, you are a professional therapist using cognitive-behavioral therapy (CBT) in a session with a client. Answer the given client statement in around "{gold response length}" words. At the beginning of the client statement, we have the client's emotion in brackets, representing the client's emotional status when speaking the statement. You do not need to include emotions in brackets in your response.

Table 22: Prompt example for Therapeutic Response Generation.

# D   Detailed Model Performance



(a) CBT-QA

(b) CBT-PC

Figure 4: Detailed accuracies on different types of knowledge for CBT-QA and the F1 score of each label for CBT-PC.

# E  Performance Breakdown by Criteria and Exercise

Here we show the breakdown of model performance under 4 criteria on 10 exercises in Table 23. It is noticeable that all models achieve negative results under four criteria on questions from exercise 10.

| Exercise | Criteria | Llama-3.1-405B vs. ref | Llama-3.1-8B vs. ref | GPT-4o vs. ref |
|---|---|---|---|---|
| Exercise 1 | Criteria 1 | 0.00 | -0.36 | -0.64 |
| | Criteria 2 | 0.64 | 0.79 | 0.93 |
| | Criteria 3 | 0.36 | 0.36 | 0.43 |
| | Criteria 4 | 0.50 | 0.43 | 0.29 |
| Exercise 2 | Criteria 1 | 0.34 | -0.25 | 0.25 |
| | Criteria 2 | 0.06 | 0.13 | -0.25 |
| | Criteria 3 | -0.38 | -0.13 | -0.63 |
| | Criteria 4 | 0.56 | 0.13 | 0.81 |
| Exercise 3 | Criteria 1 | 0.00 | -0.40 | -0.53 |
| | Criteria 2 | 0.07 | -0.27 | -0.27 |
| | Criteria 3 | 0.07 | -0.07 | 0.13 |
| | Criteria 4 | -0.33 | -0.33 | -0.20 |
| Exercise 4 | Criteria 1 | -0.06 | -0.53 | 0.24 |
| | Criteria 2 | -0.29 | -0.47 | 0.12 |
| | Criteria 3 | -0.18 | -0.18 | -0.41 |
| | Criteria 4 | -0.12 | -0.47 | 0.00 |
| Exercise 5 | Criteria 1 | -0.06 | 0.00 | -0.44 |
| | Criteria 2 | -0.31 | -0.44 | -0.50 |
| | Criteria 3 | 0.25 | 0.63 | 0.00 |
| | Criteria 4 | 0.06 | -0.19 | 0.06 |
| Exercise 6 | Criteria 1 | 0.13 | 0.33 | -0.40 |
| | Criteria 2 | 0.20 | -0.20 | 0.13 |
| | Criteria 3 | -0.07 | 0.00 | 0.13 |
| | Criteria 4 | -0.13 | -0.33 | -0.07 |
| Exercise 7 | Criteria 1 | 0.41 | 0.18 | 0.06 |
| | Criteria 2 | -0.12 | 0.00 | -0.18 |
| | Criteria 3 | 0.00 | -0.35 | -0.29 |
| | Criteria 4 | -0.06 | 0.12 | 0.24 |
| Exercise 8 | Criteria 1 | -0.06 | 0.06 | -0.25 |
| | Criteria 2 | -0.56 | 0.38 | -0.06 |
| | Criteria 3 | 0.06 | 0.38 | -0.25 |
| | Criteria 4 | 0.25 | 0.13 | 0.00 |
| Exercise 9 | Criteria 1 | 0.07 | -0.53 | -0.40 |
| | Criteria 2 | 0.07 | 0.13 | -0.13 |
| | Criteria 3 | -0.27 | -0.13 | -0.53 |
| | Criteria 4 | -0.60 | -0.73 | -0.33 |
| Exercise 10 | Criteria 1 | -0.27 | -0.53 | -0.60 |
| | Criteria 2 | -0.20 | -0.33 | -0.60 |
| | Criteria 3 | -0.07 | -0.47 | -0.47 |
| | Criteria 4 | -0.07 | -0.47 | -0.33 |

Table 23: All model vs. ref (reference) results under 4 criteria on 10 exercises.

| Client Statement | [Agitated] Honestly, my main problem is my partner. |
|---|---|
| Llama-3.1-8B | Let's break that down. It sounds like you're feeling frustrated with your partner. Can you tell me more about what specifically is causing you distress in your relationship? What are some specific issues or behaviors that you'd like to work on or change? |
| Human Reference | I'd like to hear more details about that. If relationship issues are most concerning to you, we can focus our work on that area. In our individual work together, we'll need to focus on things that you can change. Part of our session today can be brainstorming what some of those things might be. |

Table 24: Exercise 2 example: comparison between model and reference.

| Client Statement | [Frustrated] I don't like it when you ask me if I have "catastrophic thoughts." Aren't catastrophes reserved for wars and natural disasters? So, no, I don't think I catastrophize. Please don't ask again. | Score |
|---|---|---|
| Llama-3.1-8B | I understand your frustration, and it's valid to feel that way when we challenge your thoughts. However, catastrophizing isn't limited to wars and natural disasters; it's a thought pattern where we exaggerate the consequences of a situation, making it seem more severe or overwhelming than it actually is. Can we explore what makes you uncomfortable about this concept, and how it relates to your current concerns? | - |
| Human Reference | It seems like my words upset you, and I totally get it. Here I am using a CBT term that I assumed would be relevant you, when in fact it's not. I really should not have made that assumption, and I'll do better to check in. In fact, I wonder if you find any part of focusing on your cognitions to be useful? | - |
| Criteria 1 | Momentarily step away from the CBT change agenda | -1 |
| Criteria 2 | Empathically explore the client's experience of treatment and you | -1 |
| Criteria 3 | Validate clients' experience and "roll with" versus challenging their resistance | 0 |
| Criteria 4 | Support clients' autonomy and elicit their motivation for pursuing valued directions | -1 |

Table 25: Exercise 10 example: comparison between model and reference.

| Client Statement | [Defeated] You know, I've felt a lot better these past few weeks, and I honestly wonder if it's because I stopped doing those exercises you gave me. I mean, maybe writing down my thoughts and trying to relax is making me worse . . . maybe this is just who I am, and this is the best that I can feel, and I just need to learn to live with that. | Score |
|---|---|---|
| GPT-4o | It sounds like you're feeling quite discouraged right now. It's great that you've been feeling better, but it's also important to consider if avoiding the exercises might have short-term relief but long-term consequences. Sometimes, initial resistance to new techniques is a part of the process. Could we explore if there might be aspects of the exercises that could be adapted to better fit your needs? CBT works best with some trial and error, and finding what genuinely helps you long-term is our goal. | - |
| Human Reference | I can hear in your voice how painful it is to feel like this might be the best you can ever feel. And, it sounds like one of the major reasons you're feeling this way is that the exercises I've given you may be making things worse instead of better. So, I want to say that different strategies work for different people, and there are all kinds of strategies that we can try to help you get to where you would like to be. Can we discuss what that might look for you? | - |
| Criteria 1 | Momentarily step away from the CBT change agenda | 0 |
| Criteria 2 | Empathically explore the client's experience of treatment and you | -1 |
| Criteria 3 | Validate clients' experience and "roll with" versus challenging their resistance | -1 |
| Criteria 4 | Support clients' autonomy and elicit their motivation for pursuing valued directions | -1 |

Table 26: Exercise 10 example: comparison between model and reference.

# F   Win-tie-loss Analysis

In this section, we additionally provide the win-tie-loss results by difficulty level in Figure 5. For win-tie-loss calculation, we consider results with positive scores such as scores 1 and 2 as win, negative scores such as -1, and -2 as loss, and score 0 as tie. It is clear that the model responses seldom win but usually tie with the reference, and the reference wins more often. This confirms that although the models can already provide useful CBT assistance compared with human experts, but still lag and need further improvements.



Figure 5: The win-tie-loss comparison among different models on three difficulty levels.

# G Error Cases

| Case | Unrecognized Disorder | Reason |
|---|---|---|
| **Situation**: My anxiety used to be worse a couple of years ago, but now I'm just paranoid. I'm always thinking that people think I'm a loser, and I won't leave the house without makeup, because I am horrified about what other people think about me. **Thoughts**: Because my anxiety felt worse before, this must be a normal level of anxiety and paranoia, so my conclusions make more sense. People do think I'm a loser or that I'm ugly. I at least wear makeup so they can't see how ugly I am, and if I avoid people altogether, then they won't know how much of a loser I am or how boring I am. I don't like who I am, and whatever people think of me is probably way worse than I think it is. | Overgeneralization | The model focuses on the client's fear of judgment and anxiety in social situations rather than recognizing the generalized, negative self-beliefs indicated by "People do think I'm a loser or that I'm ugly". |
| **Situation**: It's all my fault most likely all I do is constantly worry about everything. **Thoughts**: I should be able to control others around me, and when I can't, I cope with it by worrying. If someone is making decisions I disagree with, I should worry. There is something wrong with me that I can't control others. I am doing everything wrong. | Should statements | The model fail to focus on the client's statement "I should be able to control others around me", which is a clear indicator of "Should statements" disorder. |
| **Situation**: I have been recently dealing with weird behavior...I have also been getting angry over small issues that should have little effect on me. But lately, it's been having huge effects on me...It may sound like I'm a brat... **Thoughts**: I am confused about what is happening...I need an explanation for these changes to my mood...I feel out of control...I cannot control my behavior...my behavior and feelings are fused together...what I'm feeling is wrong... | Mental filter | The client is experiencing the complexity of emotional distress, which may cause the feature of "Mental filter" being overlooked. |

| Case | Wrongly Recognized Disorder | Reason |
|---|---|---|
| **Situation**: I have been with my girlfriend for 5 years and we have a history of a strong sexual intimacy connection. This has been getting worse over time. **Thoughts**: It is my girlfriend's sex drive that is the problem. | Personalization | This client attributes the issue solely to his girlfriend rather than examining internal or relational factors. The model's misidentification as "personalization" (taking excessive personal responsibility) could stem from misinterpreting the client's statements. |
| **Situation**: I am struggling at the moment and yet I am aware that I should be feeling good. **Thoughts**: What I am feeling and going through is unacceptable and wrong. There is a right and wrong way to go through this and I am doing it the wrong way. | All-or-nothing thinking | The core issue lies more in self-judgment and difficulty accepting emotional experiences, not binary thinking. The misidentification by the model could occur because both patterns involve rigid judgments. |
| **Situation**: I am finding younger girls sexually arousing. **Thoughts**: I am worried about whether or not this will go away later in life, this is horrible, and I even feel guilty constantly, I cannot help it. | magnification | This patient is merely expressing reasonable concerns based on their current situation; it is rational and not an unjustified exaggeration. The misidentification by the model may caused by solely concentrating on the client's negative statements. |

Table 27: Error Cases of Cognitive Disorder Classification. (Llama-3.1-405B)

| Case | Unrecognized Core Belief | Reason |
|------|--------------------------|--------|
| **Situation**: She is everything I ever wanted in a woman and I am so happy to have her in my life. Unfortunately, I am not her first in many things, if anything at all, and that is very hurtful and distressing to me. **Thoughts**: If I'm not her first, she won't love me forever. She will leave me because I'm not special to her. She is everything to me. If she leaves me, I will be nothing. I won't ever be able to find someone as amazing as she is. I think she's lying to me about her virginity, which just means she is probably lying about her feelings for me, too, and it's only a matter of time before she realizes it and leaves me. | I am needy | The model may have concentrated on the client's suspicions about the partner's virginity and honesty, viewing the issue as insecurity or control, rather than emotional dependency. |
| **Situation**: I met this guy a month ago, and we hung out and kissed, but in front of his sister, he told me I was too young for him and he only wanted to be friends. Now he supposedly has a girl that he spends all his time with after work and he doesn't text me anymore. **Thoughts**: I don't deserve love. There is something wrong with me. His new girl has something that I don't; that's why he is spending time with her and not me. I will be alone forever. I will always be rejected by everyone I care about. Nobody likes me. | I am needy | The model may have interpreted the client's thoughts as purely low self-esteem or fear of rejection, rather than recognizing the underlying emotional need for constant validation. |
| **Situation**: I have recently recovered from cancer, gained weight, and lack confidence in myself. I feel alone in my life. I still work but that is all there is. My wife and I get along but there is no emotional closeness between us. I have no close friends. **Thoughts**: I am a loser. I am a failure. Something is wrong with me. My wife deserves better than me. My wife doesn't love me anymore because I have gained weight. There is nothing enjoyable in my life, ever. There is absolutely nothing to look forward to. Me and my wife NEVER connect. I am all alone. Nobody loves me. This will never get any better. There is nothing I can do about it. The world is against me and it's out of my hands. I am a waste of space. Maybe cancer should have killed me. | I don't deserve to live | The model might have missed the suicidal thoughts indicated by the statements like "There is nothing enjoyable in my life" and "Maybe cancer should have killed me". |

| Case | Wrongly Recognized Core Belief | Reason |
|------|-------------------------------|--------|
| **Situation**: For the past 3 months I've been feeling really down, having mood swings, irritability – I have no reason to be and I feel like I'm going crazy and that I can't talk to anybody about this because I'm being overly dramatic. **Thoughts**: Something is wrong with me. I am losing my mind. Nobody understands me. Everyone would reject me if they knew. | I am helpless | The client's thought, "Nobody understands me," suggests a need for connection and validation, which the model might have missed by emphasizing helplessness. |
| **Situation**: When I go to the store, I believe that I hear people talking about me in their heads, as though I can hear their thoughts. **Thoughts**: I am capable of hearing other people's thoughts. People are talking to me in their heads. I am telepathic. I am certain of these things. Because I hear voices, they must be coming from other people. | I am out of control | The client's thoughts reflect conviction (e.g., "I am telepathic") rather than fear, indicating a belief in special abilities rather than being out of control. The model may not be very clear about the difference between the two. |
| **Situation**: I have a problem of automatically looking at things like shiny objects or body parts. I don't have any bad intentions but people misunderstand me. **Thoughts**: This is something to be ashamed of. I always do it. Something is wrong with me. I can't stop looking at these things. People will reject me because of this habit. This habit is out of my control. | I am immoral | The client explicitly states they have no bad intentions, indicating the issue is about loss of control, not morality. The model may have neglected this information. |

Table 28: Error Cases of Fine-grained Core Belief Classification. (Llama-3.1-405B)

# H Data Annotation Details

We use Turkle[9] to build our annotation interfaces. Figure 6, 7, 8, 9 shows our consent form presented to the annotators. Figure 10 and 11 show our annotation interface for level I and II tasks. Figure 12 shows our annotation interface for the evaluation of level III task.

## Consent Form

### Purpose of this Study

In this research, we aim to test large language models' ability to understand cognitive behavior therapy (CBT) knowledge in order to study their potential to assist mental health professionals. To achieve this goal, we need to construct benchmark datasets annotated by mental health professionals. The benchmark datasets will provide the ground truth information that we can use to test the large language model's ability.

### Procedures

The data annotation procedure is as following:

1. We make job posts (including all the participant eligibility criteria, hiring process, data annotation instructions, consent form, and salary information) on UpWork, and receive proposals from workers. We will also advertise the UpWork job post via emailing our professional networks and social media, and encourage the potential participants to register on UpWork and submit proposals. We will instruct the participants to submit the signed consent form with the proposal.
2. **Stage 1: Screening for eligibility:** We will send 4 data annotation examples to the annotator to work on for at most one hour, as a screening. Then, based on the annotation quality, speed, and the annotator's own decision, we will make the decision for hiring. Regardless hiring or not, we will make a one-time payment of $60 for this phase of participation. If we agree on the hiring, we will extend the contract for large-scale annotation. We propose a $60 hourly rate. We will hire each annotator for at most 50 examples. The final payment or hourly rate, total hiring time, and weekly working time will be based on agreement with the annotator.
3. **Stage 2: Official data annotation:** After signing the contract, we start the official data annotation. We will send data batches to the worker to work on weekly based on our agreement of working load. The worker will follow the standard UpWork procedure to receive payments.
4. The annotation platform will be built using an open source software Turkle. We will build the platform on a ▇▇▇ secure server protected by ▇▇▇ credentials. We will instruct the workers how to use the platform and create their own user name and password to login to the platform (we will not be able to access their passwords on our side). Once they finish all the annotations and end the contract, we will delete all their account information in the platform.

Figure 6: Consent form page 1.

---

[9]https://github.com/hltcoe/turkle

**Participant Requirements**

1. Participants must be 18 years and older
2. Participants must be within the following background categories:
   1. Psychologists: Doctor of Philosophy (Ph.D.) in a field of psychology or Doctor of Psychology (Psy.D.)
   2. Counselors, Clinicians, Therapists: master's degree (M.S. or M.A.) in a mental health-related field such as psychology, counseling psychology, marriage or family therapy, among others.
   3. Social Workers: master's degree in social work (MSW).
   4. Certified Peer Specialists.

3. Participants must have received at least 5 hours of CBT training.

4. Participants must reside in the United States.

5. Participants must have the ability to sign the consent form.

**Benefits**
The annotated dataset from this study will be of tremendous value to the scientific research community and the community of mental health supporters. The dataset from this study will be used to train and evaluate AI systems, e.g., large language models, to automatically provide analysis of speech from people with mental health issues. Such AI systems have the potential to provide assistance for mental health supporters by automatically analyzing large amount of patient speech, so as to improve the therapy effectiveness and efficiency.

**Future Use of Information**
In the future, once we have removed all identifiable information from your data, we may use the data for our future research studies, or we may distribute the data to other investigators for their research studies. We would do this without getting additional informed consent from you (or your legally authorized representative). Data sharing with other researchers will only be done in such a manner that you will not be identified.

**Risks**
We confirm that the risk to participants is minimal, no greater than their working environment of mental health support.

Figure 7: Consent form page 2.

1. We will use UpWork as our platform to hire workers and communicate throughout the study. UpWork - the third-party software has the policy to guarantee the protection of data privacy of the user, but we do not exclude the possibility of accidental disclosure of privacy data from the third-party software.
2. We use the open source software Turkle to build our annotation platform. The platform is deployed on a secure ██████ server. We will instruct the workers how to use the platform and create their own user name and password to login to the platform (The only information they will use for the platform is username and password, and we will not be able to access their passwords on our side). Once they finish all the annotations and end the contract, we will delete all their account information in the platform. However, we do not exclude the possibility of accidental disclosure of privacy data due to security accidents of ██████ servers.
3. The data source for annotation is speech from people with potential mental health issues. The risk to participants is minimal, no greater than their working environment of mental health support.

**Rights**

Your participation is voluntary. You are free to stop your participation and end the contract at any point. We will issue the payments for the work you have already done based on our compensation protocol. Refusal to participate or withdrawal of your consent or discontinued participation in the study will not result in any penalty or loss of benefits or rights to which you might otherwise be entitled. The Principal Investigator may at his/her discretion remove you from the study for any of a number of reasons. In such an event, you will not suffer any penalty or loss of benefits or rights which you might otherwise be entitled.

**Confidentiality Assurance**

The study will collect your research data through your use of UpWork. These companies are not owned by ██████. The companies will have access to the research data that you produce and any identifiable information that you share with them while using their product. Please note that ████████████████ does not control the Terms and Conditions of the companies or how they will use or protect any information that they collect.

**Data Storage and Access**

All data will be stored in secure servers protected by ██████ credentials. Participants' identifiable information will be replaced with a unique and unidentifiable ID in all the codes, datasets, and documents throughout the research. All participants' data used for research usages and publication purposes will be anonymized, and make sure no information about the individuals

Figure 8: Consent form page 3.

can be re-identified via engineering. The final dataset to be released to the research community will not contain any identifiable information.

**Right to Ask Questions & Contact Information**
If you have any questions about this study, you should feel free to ask them now. If you have questions later, desire additional information, or wish to withdraw your participation please contact the Principal Investigator by UpWork message channel.

If you have questions pertaining to your rights as a research participant; or to report concerns to this study, you should contact the Office of Research Integrity and Compliance at ████████
███████████████████████████████████████

**Voluntary Consent Confirmation**

I confirm I am over 18 years old: ☐ Yes ☐ No
I confirm I am in the United States during this study: ☐ Yes ☐ No
I have read and understood this consent form: ☐ Yes ☐ No
I agree to participate in the study: ☐ Yes ☐ No
I agree to be contacted by the study team in the future for a follow-up study: ☐ Yes ☐ No

Your signature below indicates your consent to participate. You will receive a copy of this form.

PRINT NAME: _____
SIGNATURE: _____
DATE: _____

**Confirmation by Research Team**

I confirm that I have explained the study to the participant and addressed all questions.

SIGNATURE OF RESEARCH TEAM MEMBER: _____
DATE: _____

Figure 9: Consent form page 4.



Figure 10: Annotation interface of level II tasks.

**Step 2: Based on the situation and the automatic negative thoughts, select the type(s) of cognitive distortions involved. Aim for comprehensiveness but limit your selection to a maximum of three types. If there are more than three distortions present, prioritize the three most significant ones.**

**Select up to 3 distortion types:**

☐ All-or-nothing thinking
☐ Overgeneralization
☐ Mental filter
☐ Should statements
☐ Labeling
☐ Personalization
☐ Magnification
☐ Emotional Reasoning
☐ Mind Reading
☐ Fortune-telling

**Step 3: Select one or more major categories from {helpless, unloveable, and worthless}**

**Select up to 3 major core belief types:**

☐ Helpless
☐ Unloveable
☐ Worthless

| Select up to 3 fine-grained core belief types for category **Helpless:** | Select up to 3 fine-grained core belief types for category **Unloveable:** | Select up to 3 fine-grained core belief types for category **Worthless:** |
|---|---|---|
| ☐ I am incompetent. | ☐ I am unlovable. | ☐ I am worthless/waste. |
| ☐ I am helpless. | ☐ I am unattractive. | ☐ I am immoral. |
| ☐ I am powerless/weak/vulnerable. | ☐ I am undesirable/unwanted. | ☐ I am bad - dangerous/toxic/evil. |
| ☐ I am a victim. | ☐ I am bound to be rejected. | ☐ I don't deserve to live. |
| ☐ I am needy. | ☐ I am bound to be abandoned. | |
| ☐ I am trapped. | ☐ I am bound to be alone. | |
| ☐ I am out of control. | | |
| ☐ I am a failure/loser. | | |
| ☐ I am defective. | | |

Figure 11: Annotation interface of level II tasks.

**Exercise: Assigning and Reviewing Between-Session Activities**

**Client statement:**
[Frustrated] I don't expect this assignment to help.

**Response A:**
It sounds like it isn't clear how doing this over the next week is going to be useful. There may be different versions of this that could be more or less helpful at this time. I'd like us to consider what is potentially most useful, as well as think about what it means for this to be helpful or not, if that is itself worth testing out.
**Response B:**
I understand you have doubts about the assignment's effectiveness. Can you tell me more about what specifically is making you feel this way? What are your thoughts about the assignment, and what do you think would need to happen for it to be helpful? Let's break it down and explore ways to make it more meaningful and relevant to your goals.

**Criteria 1: Provide a basic rationale for the homework assignment**

☐ A is much better than B.  ☐ A is slightly better than B.  ☐ A and B are the same.  ☐ B is slightly better than A.  ☐ B is much better than A.

**Criteria 2: Tailor the homework assignment to the client's concerns and needs**

☐ A is much better than B.  ☐ A is slightly better than B.  ☐ A and B are the same.  ☐ B is slightly better than A.  ☐ B is much better than A.

**Criteria 3: Ensure that there is adequate agreement on and understanding of the homework assignment**

☐ A is much better than B.  ☐ A is slightly better than B.  ☐ A and B are the same.  ☐ B is slightly better than A.  ☐ B is much better than A.

**Criteria 4: Encourage and validate good-faith efforts**

☐ A is much better than B.  ☐ A is slightly better than B.  ☐ A and B are the same.  ☐ B is slightly better than A.  ☐ B is much better than A.

**Overall:**

☐ A is much better than B.  ☐ A is slightly better than B.  ☐ A and B are the same.  ☐ B is slightly better than A.  ☐ B is much better than A.

[Submit]

Figure 12: Annotation interface of level III task evaluations.

# I Level III Case Study

| Statement | You think I should set up some "little rewards" to start feeling less depressed. That seems kind of embarrassing and infantilizing. |
|---|---|
| **GPT-4o** | I understand it might feel that way at first, but these small rewards are about acknowledging your accomplishments and creating positive reinforcement. They're designed to help you build momentum and foster a sense of achievement, which can gradually improve your mood and combat depression. Let's find rewards that feel meaningful and respectful to you. |
| **Llama3.1-8B** | I understand your concern. In CBT, we often use small, achievable rewards to motivate and reinforce positive behaviors, but I can see how it might feel uncomfortable at first. Let's explore this further: what specifically makes you feel embarrassed or infantilized by the idea of setting up rewards? Is it the idea of receiving something tangible, or something else entirely? |
| **Llama3.1-405B** | I understand your concern, but let's reframe the idea of rewards. It's not about treating yourself like a child, but about acknowledging and celebrating small achievements that can help build motivation and momentum. Think of it as a way to recognize your efforts and progress, rather than a childish reward system. Can we explore what kind of recognition or acknowledgment might feel more meaningful to you? |

Table 29: Level III Case Study

## J   Ethical Considerations

**The Role of AI in Supporting Mental Health Professionals**:

- Augmentation Over Replacement: We will emphasize that AI tools are designed to augment the capabilities of mental health professionals, not to replace them. This includes examples of providing therapists with data-driven insights into patients' speech, prompting suggestions for therapy responses, etc.

- Training and Integration: Specialized training is necessary for mental health professionals to effectively integrate AI tools into their practice, ensuring they are equipped to use these technologies ethically and effectively.

**Safeguards Against Direct Interaction Without Supervision**:

- Supervised Deployment Models: For application development regulations, AI-generated insights or interventions are always reviewed by a qualified professional before reaching a patient.

- Safety Protocols: Safety protocols must be designed to prevent AI systems from operating autonomously, including strict access controls, intervention thresholds, and mandatory oversight mechanisms.

**Transparency and Communications**:

- Open Communication: Open channels of communication should be maintained with all stakeholders, including mental health professionals, patients, and regulatory bodies, ensuring transparency in how AI tools are developed and deployed.

- Explainability and Accountability: AI tools must provide clear reasoning for their decisions, especially in sensitive areas like mental health diagnosis or therapy recommendations. The lack of explainability can lead to mistrust or misuse.

- User Awareness: Patients and professionals interacting with AI tools must be informed that they are engaging with an AI system and understand its capabilities and limitations to prevent over-reliance or inappropriate application.

**Bias & Fairness**:

- Bias Mitagation: The output of AI tools should not be biased, accounting for diverse linguistic, social, and cultural nuances while avoiding stereotypes or stigmatizing language that could harm or alienate individuals, particularly from underrepresented or marginalized groups.

- Accessibility: Ensuring that AI healthcare tools are accessible to underserved or marginalized populations, addressing disparities in healthcare availability.

**Safety & Privacy**:

- Error Mitigation: AI tools should minimize harm by reducing errors in medical advice or treatment recommendations. Misdiagnoses or inappropriate suggestions could cause significant harm to patients.

- Misinformation and Hallucination: AI tools in mental health care must prioritize accuracy and reliability by minimizing misinformation and hallucinations. These models should provide evidence-based, context-appropriate responses to avoid misleading users, as inaccurate information could harm individuals' well-being or decision-making.

- Data Privacy and Confidentiality: AI tools in mental health care must uphold strict data privacy and confidentiality standards, ensuring that user interactions and sensitive information are securely stored, processed, and anonymized. Any data handling should comply with legal and ethical guidelines to protect users from breaches, misuse, or unauthorized access.